

# Two Sample Studies

(+ a few first words on the user study “mindset”)

**Ulrike Kuhl**

*Bielefeld University, Bielefeld, Germany*

# Designing a User Study

To what end do you want to test what with whom in which context?

# Designing a User Study

To what end do you want to test what with whom in which context?

Purpose/Objective: Examining...



...humans



...models



...human-AI interaction



...explanations

# Designing a User Study

To what end do you want to test what with whom in which context?

Purpose/Objective: Examining...



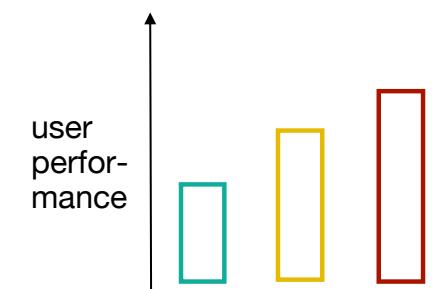
Focus/Variables - examples:

If the odor had been less foul, the food would have been deemed edible.

vs. Lime ([Ribeiro et al., 2016, ACM SIGKDD](#))



--> compare XAI methods with each other



--> assess performance, subjective measures, etc.

# Designing a User Study

To what end do you want to test what with whom in which context?

Target User Group:



VS.

VS.

--> Challenges: sample size, recruitment, inclusion / exclusion criteria

# Designing a User Study

To what end do you want to test what with whom in which **context**?

Target User Group:



VS.



VS.



--> sample size, recruitment,  
inclusion / exclusion criteria

Application / Environment:

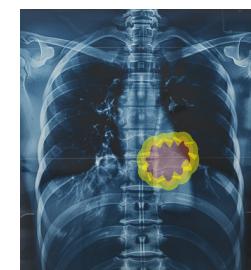


VS.



VS.

„Imagine  
you're a  
gardener...“



Med-AI detected cancer.  
Do you agree?

--> User's world knowledge matters!

# What has been done: Sample Study I

To what end do you want to test what with whom in which context?

To what end?

Examine impact of continuous and categorical features in counterfactual and causal explanations

"if only..."

"because..."

James	
Gender	Male
Weight	81kg

categorical

continuous

# What has been done: Sample Study I

To what end do you want to test what with whom in which context?

To what end?

Examine impact of continuous and categorical features in counterfactual and causal explanations

"if only..."

"because..."

James

Gender	Male
Weight	81kg

categorical

continuous

What?

Human prediction accuracy w.r.t. a system's decisions + subjective satisfaction & trust in the system

# What has been done: Sample Study I

To what end do you want to test what with whom in which context?

Whom?

“Representative”  
human users



# What has been done: Sample Study I

To what end do you want to test what with whom in which context?

Whom?

“Representative”  
human users

Context?

*SafeLimit* App

James	
Gender	Male
Weight	81kg
Units	6
Duration	105 mins
Stomach	Full
Limit	?

Over the limit



Don't know



Under the limit



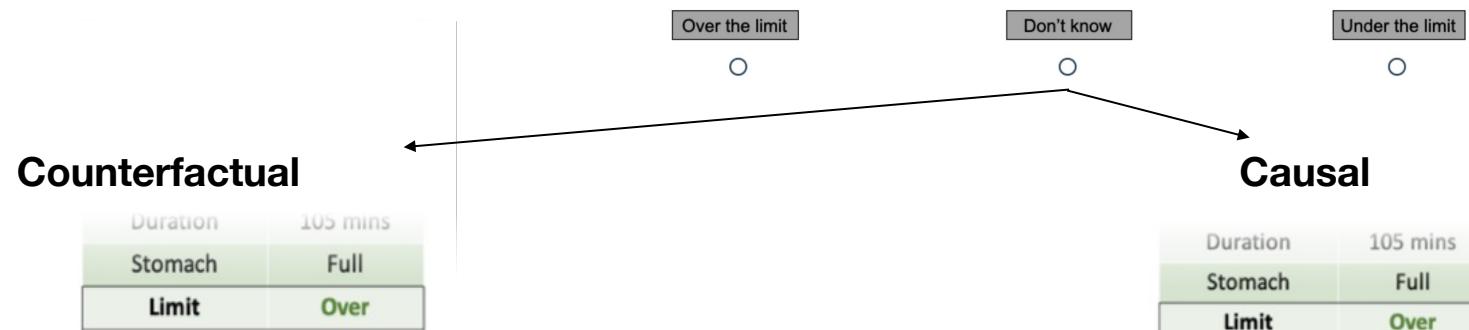
# What has been done: Sample Study I

To what end do you want to test what with whom in which context?

**Whom?** “Representative” human users

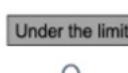
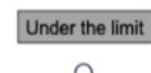
James	
Gender	Male
Weight	81kg
Units	6
Duration	105 mins
Stomach	Full
Limit	?

**Context?** SafeLimit App



**Explanation**  
If James had drunk 5 units instead of 6 units, he would have been under the limit.

**Explanation**  
James is over the limit because he drank 6 units.



# What has been done: Sample Study I

## Between-Subjects

--> 3 groups in 3 conditions

Group 1: Counterfactual



vs.

Group 2: Causal



vs.

Group 3: Control



### Explanation

If James had drunk 5 units instead of 6 units, he would have been under the limit.

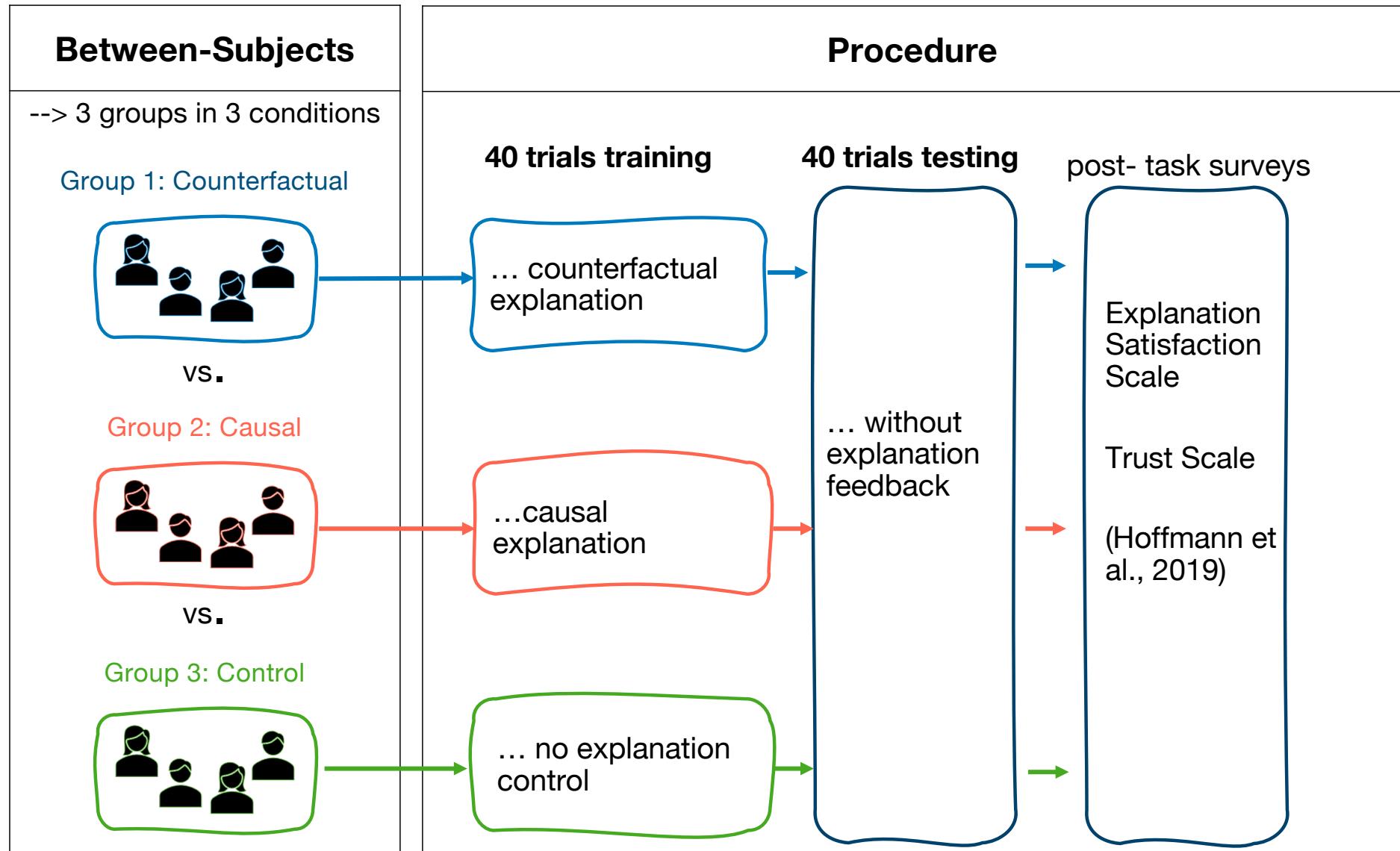
### Explanation

James is over the limit because he drank 6 units.

### Explanation

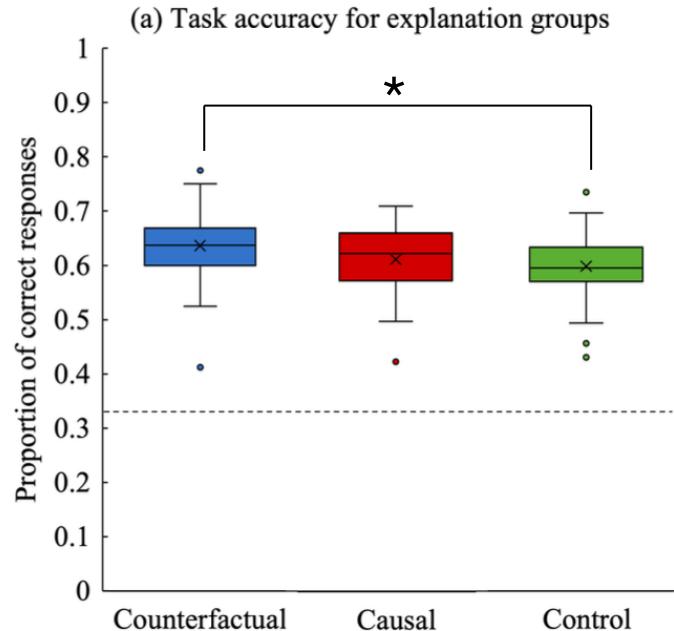
James is over the limit.

# What has been done: Sample Study I

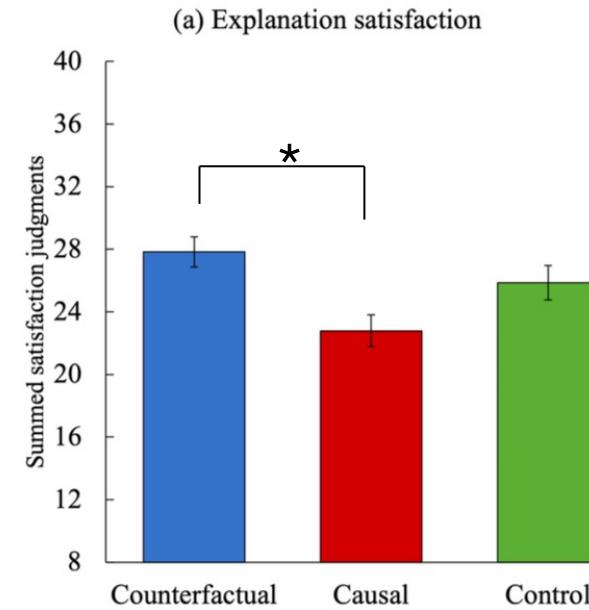


# What has been done: Sample Study I

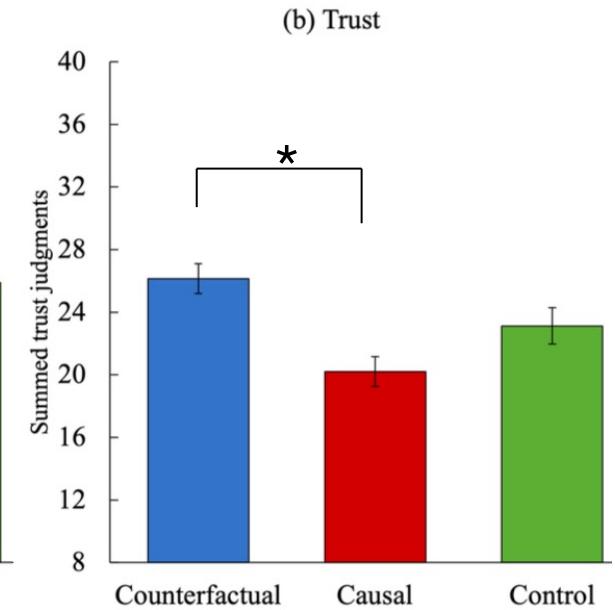
How accurate?



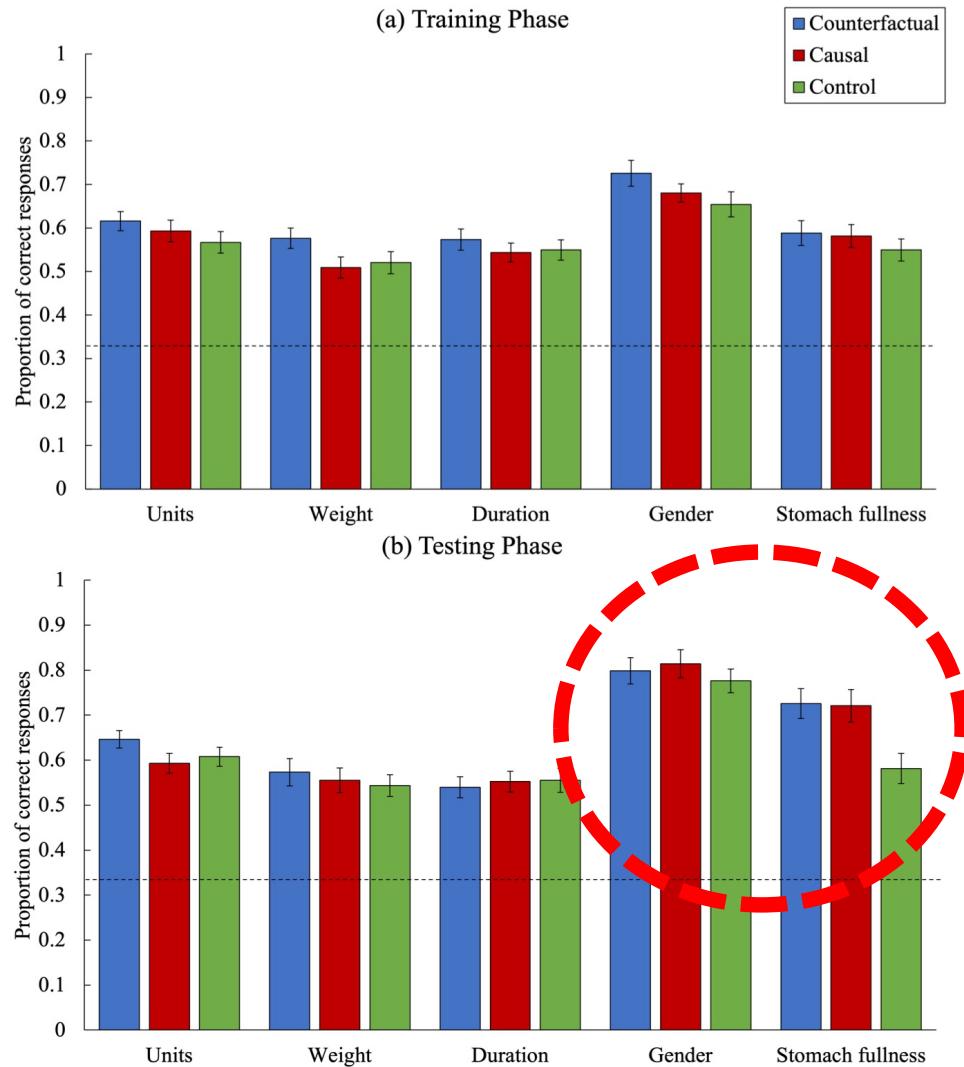
How satisfied?



How trusting?



# What has been done: Sample Study I



→ accuracy improved from training to test for categorical, but not for continuous features

# What has been done: Sample Study I

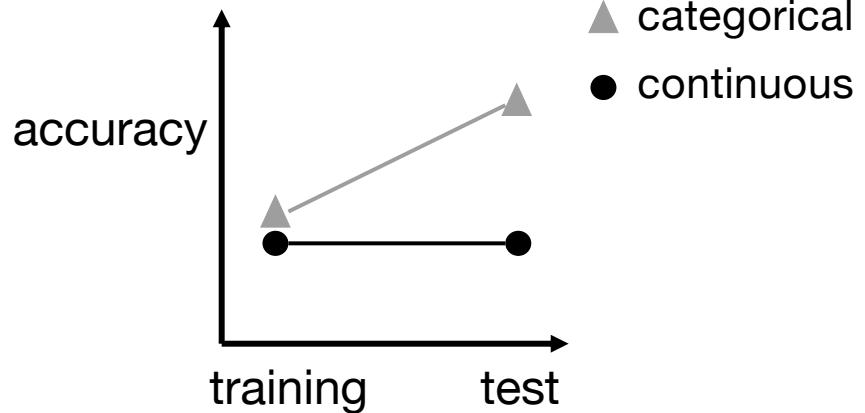
## To what end?

Examine impact of continuous and categorical features in counterfactual and causal explanations

### Accuracy Improvement:

- generally:  
explanations > no explanations
- specifically:  
CF > no explanations

### Feature Types:



### Subjective Evaluations (satisfaction and trust):

- CF > causal

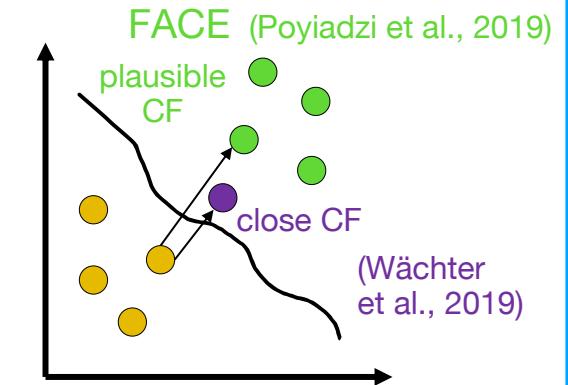
→ explanation type and fundamental aspects of features matter!

# What has been done: Sample Study II

To what end do you want to test what with whom in which context?

To what end?

Examine usefulness of a plausibility constraint for generation of counterfactual explanations

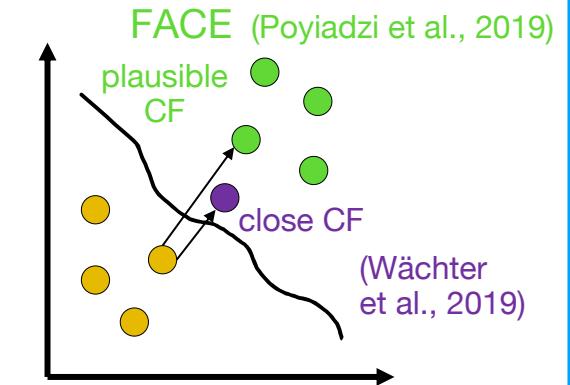


# What has been done: Sample Study II

To what end do you want to test what with whom in which context?

To what end?

Examine usefulness of a plausibility constraint for generation of counterfactual explanations



What?

Human learning of an unknown system, supported by explanations

# What has been done: Sample Study II

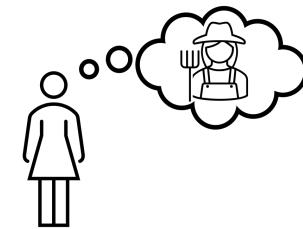
To what end do you want to test what with whom in which context?

**Whom?**

“Representative”  
human users

**Context?**

*Abstract scenario:  
Alien Zoo framework*



... selects plants...



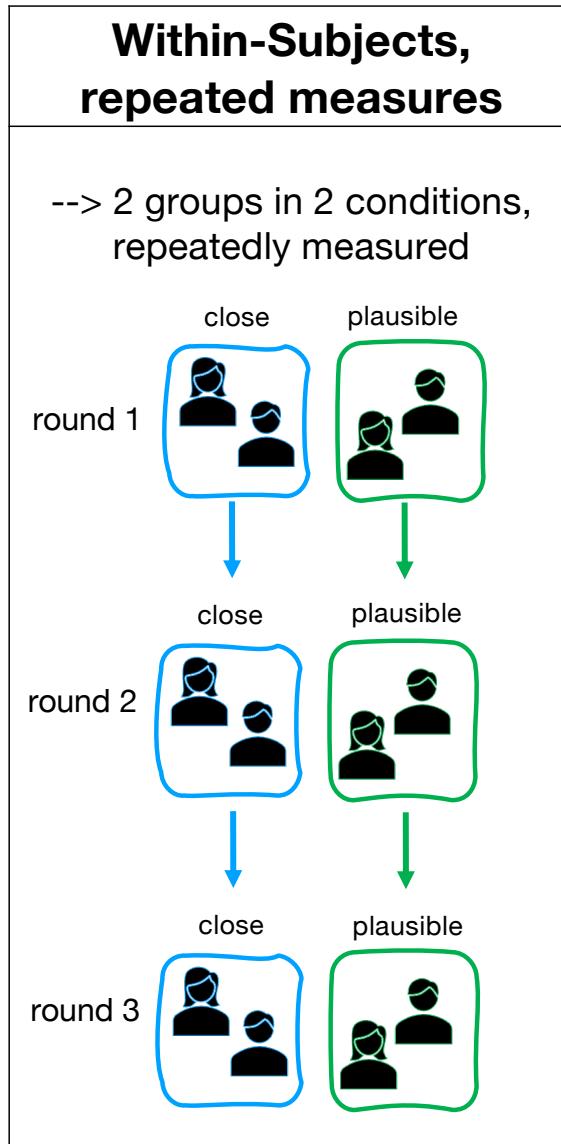
... to feed aliens:



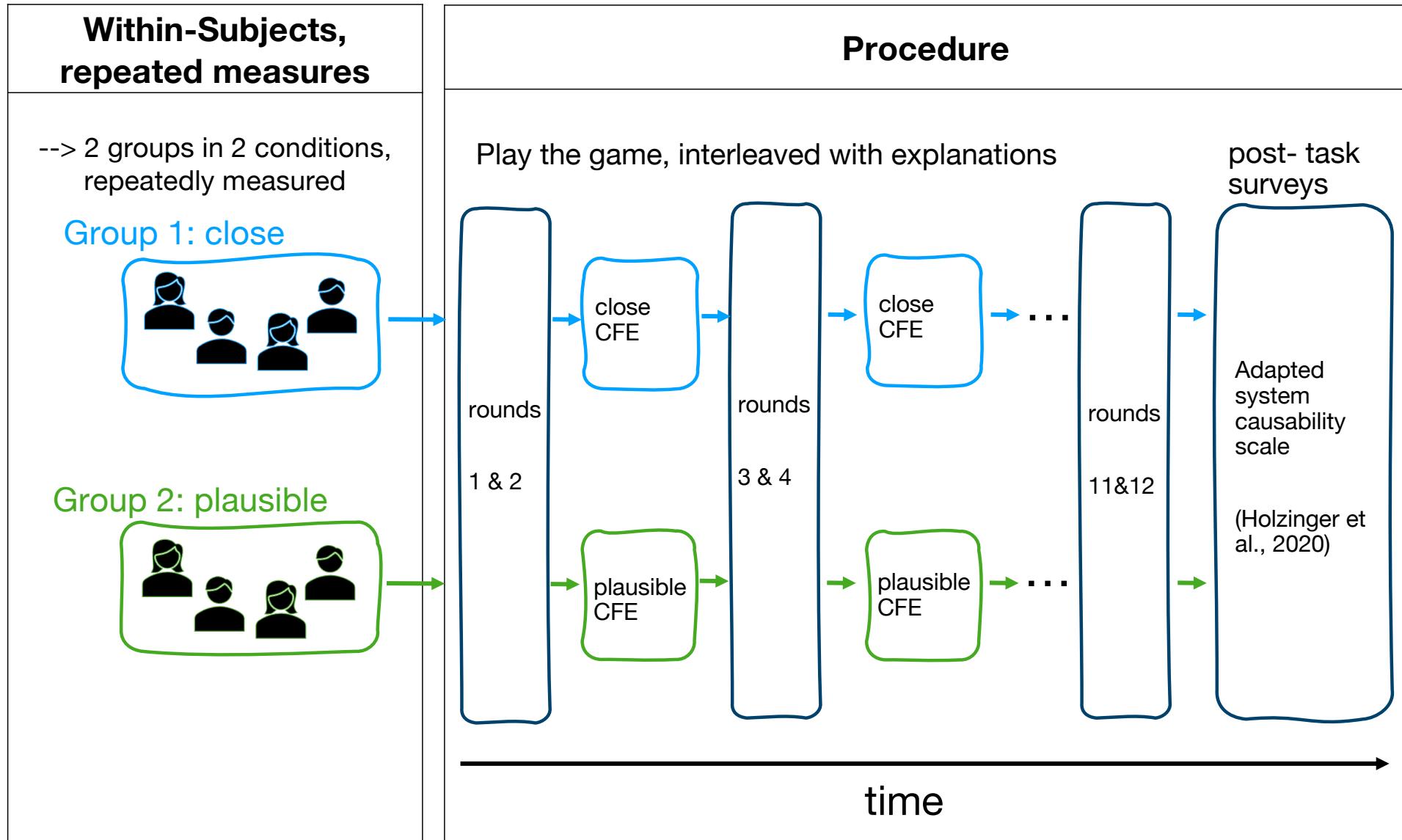
Good choice



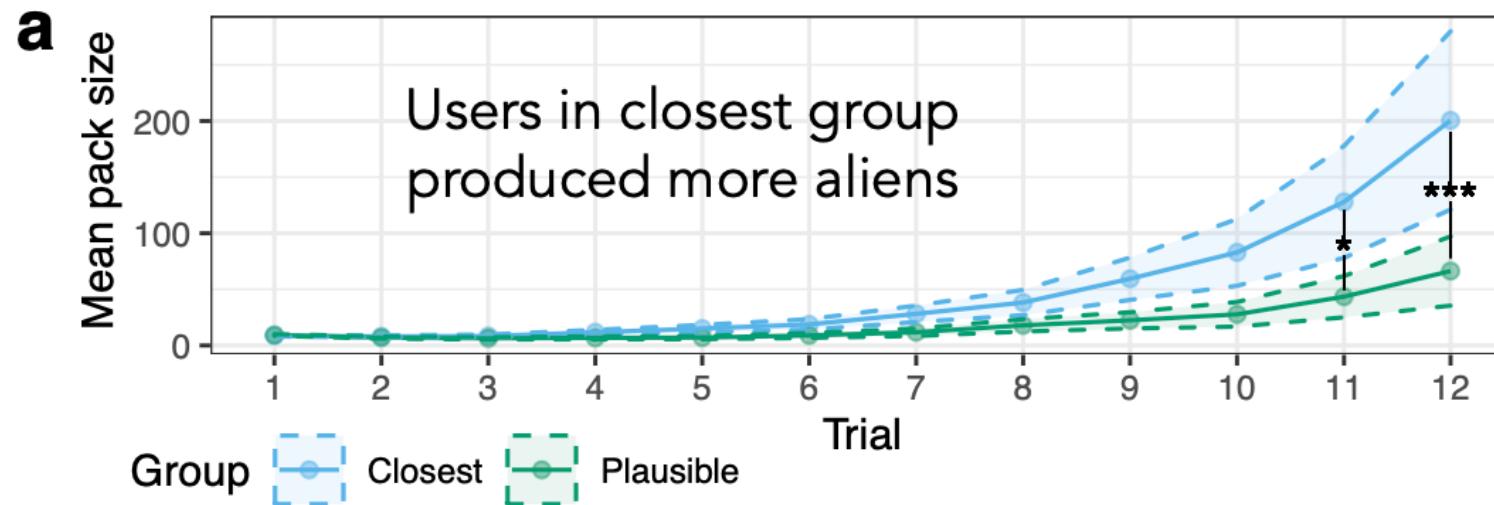
# What has been done: Sample Study II



# What has been done: Sample Study II



# What has been done: Sample Study II



Yet:

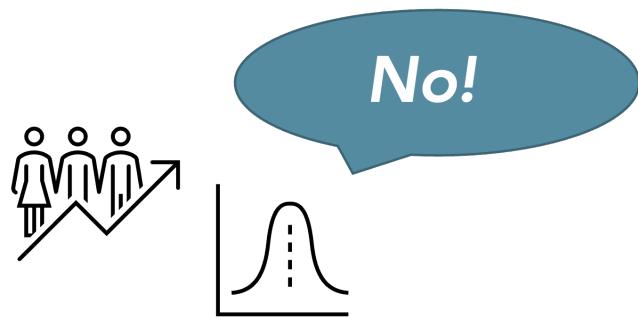
- participants in both groups could distinguish important from unimportant features
- no differences in subjective measures were detected

# What has been done: Sample Study II

To what end?

Examine usefulness of a plausibility constraint for generation of counterfactual explanations

Do **novice users** benefit from a plausibility constraint on computing CFEs, when tasked to **gain new information from an unknown system in an abstract domain?**



No!

Staying more similar to the original instance is more effective for learning, if no prior knowledge exists!

# Plan for the Day

CF Tutorial	
TIME	Topics
9:00 AM	<b>Introduction</b>  Hello and Introducing Ourselves!  <b>Hands-on:</b> Trying Our Study (follow link)
9:30 AM	<b>Historical Fundamentals of Counterfactuals</b>  From Philosophy to XAI (via Psychology)  Two Sample User Studies and Q&A
10:30 AM	COFFEE (10:30-11:00)
11:00 AM	<b>Fundamentals of Counterfactuals in AI</b>  Formalisation  Modelling Approaches & Key Constraints
11:30 AM	<b>Using Counterfactual Algorithms</b>  <b>Hands-on:</b> A Counterfactual Toolbox (AA)  <b>Hands-on:</b> Checking Out Notebooks and Q&A
12:00 PM	<b>Fundamentals of User Studies</b>  User Studies I: A Simple Two-Group Design
12:30 PM	LUNCH (12:30-14:00)
2:00 PM	<b>Algorithmic Growth Points</b>  Computational Future Directions and Q&A
2:30 PM	<b>More Fundamentals of User Studies</b>  User Studies I: A Simple Two-Group Design (cont.)
3:00 PM	COFFEE (15:00-15:30)
3:30 PM	<b>From Fundamentals to an Actual User Study</b>  User Studies II: A More Complex Design  User Studies III: Even More Complex Designs  <b>Hands-on:</b> Looking At Our Study
5:00 PM	<b>Closing Session, Discussion and Final Q&amp;A</b>
	TUTORIAL END

